

Extending BDI Plan Selection to Incorporate Learning from Experience

Dhirendra Singh & Sebastian Sardina & Lin Padgham

RMIT University, Melbourne, Australia

Abstract

An important drawback to the popular Belief, Desire, and Intentions (BDI) paradigm is that such systems include no element of learning from experience. We describe a novel BDI execution framework that models context conditions as *decision trees*, rather than boolean formulae, allowing agents to *learn* the probability of success for plans based on experience. By using a probabilistic plan selection function, the agents can balance exploration and exploitation of their plans. We extend earlier work to include both parameterised goals and recursion and modify our previous approach to decision tree confidence to include large and even non-finite domains that arise from such consideration. Our evaluation on a pre-existing program that relies heavily on recursion and parametrised goals confirms previous results that naive learning fails in some circumstances, and demonstrates that the improved approach learns relatively well.

Keywords:

BDI Intelligent Agents, Machine Learning, Hybrid Systems

1. Introduction

Agents are an important technology that have the potential to take over contemporary methods for analysing, designing, and implementing complex software systems suitable for domains such as telecommunications, industrial control, business process management, transportation, logistics, and aeronautics [1, 2, 3, 4]. The BDI model of agency [5, 6] is a popular and well-studied

Email address: {firstname.lastname}@rmit.edu.au (Dhirendra Singh & Sebastian Sardina & Lin Padgham)

approach with substantial theoretical and practical work. It has its roots in philosophy with Bratman’s [7] theory of practical reasoning and Dennett’s theory of intentional systems [8]. A recent industry study [9] analysing several applications claimed that the use of BDI (Belief-Desire-Intention) agent technology in complex business settings can improve overall project productivity by up to 500%. Also the agent approach allowed the business to change and extend solutions quickly helping to bridge the semantic gap between the business side and IT development. BDI systems have built into them an ability to balance pro-actively pursuing a goal, with reactively responding to the environment. They also have a well developed failure recovery mechanism. This makes them very suitable for robotics applications operating in a physical world which is often more error prone than a software domain.

BDI systems, despite their strengths, do not however incorporate any ability to learn from experience. Our work makes a start at addressing this issue, focussing specifically on learning which plan to select next, to resolve a particular goal in a particular world state.

There are many agent programming languages and development platforms in the BDI tradition, including JACK [10], JADEX [11], and Jason [12] among others. All of them follow a similar basic architecture, whereby *abstract plans* written by programmers are combined and used *reactively* in real-time, in a way that is both flexible and robust. Concretely, a BDI agent is built around a *plan library*, a collection of pre-defined *hierarchical plans* indexed by goals and representing the standard operational procedures of the domain (e.g., landing a plane). A *context condition* attached to each plan states the conditions under which the plan is a sensible strategy to address the corresponding goal in a given situation (e.g., it is not raining). The execution of a BDI system then relies on *context sensitive subgoal expansion*, allowing agents to “act as they go” by making *plan choices* at each level of abstraction with respect to the current situation. Although this is quite flexible and effective, an important drawback is the lack of ability to learn from ongoing experience. An ability to *learn* plan selection in particular situations, adds a whole new layer of robustness and flexibility. Firstly there may be situations where it is difficult to determine in advance the exact situation under which a particular approach is likely to succeed. This is especially the case when it involves complex combinations of values of environmental variables. Secondly, an environment may change over time, or be slightly different in different deployment locations. The ability for the agent system to observe and learn from its performance is obviously a very desirable property.

In our work we achieve this by replacing (or augmenting) the usual

boolean formula for representation of context conditions, by a decision tree [13] which is learnt based on experience. Our plan selection is then based on a probabilistic approach, usually choosing the plan which has the highest likelihood of success, based on experience. This probabilistic approach is also more suitable than the standard boolean approach for complex and often partially observable worlds, where various plans may be worth trying, but have different chances of success.

There are however various nuances that must be addressed for such *online* learning. There is the usual balance between exploration and exploitation evident in all learning. Moreover, learning is impacted by the structure imposed by the hierarchical representation of BDI programs. We have addressed these issues in previous papers [14, 15], looking at various approaches to the problem. In this paper we add the ability to deal with parameterised goals and with recursive calls, both of which are essential for real applications. Unfortunately, once we add this expressivity our previous preferred approach does not scale. Consequently we develop a simplified approximation to achieve the same basic intuition which we have previously shown to be correct in principle. We then empirically evaluate our approach by taking an existing BDI program from the JACK tutorial, removing the context conditions, and learning the appropriate use of the plans provided using our framework.

Our approach can easily be combined with the standard plan selection mechanism, by allowing the agent programmer to provide initial context conditions that could later be automatically “refined” by the agent system. By doing so, one can effectively take a BDI program and “tune it” using our learning framework. For simplicity, though, context conditions are learnt from scratch in our experimental work.

In the next section we introduce both BDI programming and our learning framework, as well as an overview of our previous approaches. We then describe in detail the learning framework that incorporates the additional aspects of parameterised goals and recursive calls, with our revised approach to address previously identified issues. We show empirical evaluation on an example program developed by Agent Oriented Software for their JACK tutorial, by removing the context conditions and applying our learning approach. We finish with a discussion of outstanding issues and related work.

2. Preliminaries

2.1. BDI Agent Systems

BDI agent-oriented programming is a popular, well-studied, and practical paradigm for building intelligent agents situated in complex and dynamic en-

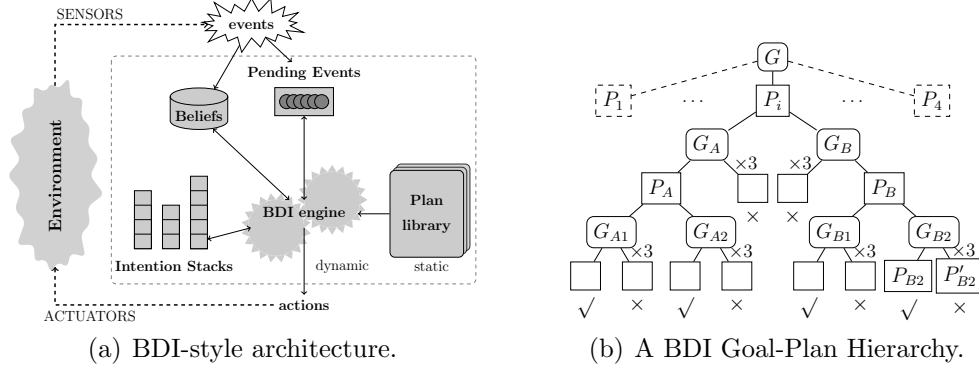


Figure 1: BDI Architecture and Goal-Plan Hierarchy.

vironments with (soft) real-time reasoning and control requirements [16, 9]. Generally speaking, BDI agent-oriented programming languages are built around an explicit representation of propositional attitudes (e.g., beliefs, desires, intentions, etc.). A BDI architecture addresses how these components are represented, updated, and processed to determine the agent’s actions.

A BDI agent consists, basically, of a belief base (akin to a database) which stores the agent’s knowledge about the world, a set of pending event-goals,¹ which include both external percepts or messages and internal goals, a plan library, and an intention structure. Figure 1(a) depicts a typical BDI architecture. The plan library contains rules of the form $e : \psi \leftarrow \delta$ indicating that δ is a suitable procedure for achieving event-goal e when context condition ψ is true. Among other operations, the plan body procedure δ will typically include the execution of actions (*act*) in the environment and subgoals (! e) that are in turn resolved by selecting suitable plans for those subgoals. For example, the following plan rules may be part of the plan library of an elevator controller:

$$\begin{aligned} \text{Serve}(\text{floor}) : \quad & \text{Serving}(\text{floor}) \leftarrow !\text{GoTo}(\text{floor}); \text{Open}; \text{Close}; \text{Off}(\text{floor}) \\ \text{GoTo}(\text{floor}) : \quad & \text{At}(x) \wedge x > \text{floor} \leftarrow \text{GoUp}; !\text{GoTo}(\text{floor}) \end{aligned}$$

That is, to serve a request from a floor (i.e., event $\text{Serve}(\text{floor})$) which the elevator is supposed to serve (i.e., condition $\text{Serving}(\text{floor})$ is true), can be

¹In this paper the terms *event-goal*, *event*, and *goal* are used interchangeably.

achieved by first going to the floor (by solving the $!GoTo(floor)$ sub-goal), then opening and closing the door, and finally turning the floor’s request light off. In turn, to go to a floor that is above the current location of the elevator, it needs to go up one floor (i.e., execute primitive action $GoUp$) and then post again the (sub)goal of reaching the floor in question.

The basic *reactive goal-oriented behavior* of BDI systems involves the system responding to events by selecting an appropriate plan from the library, and placing its program body into the intention base, a structure containing the current, partially instantiated, plans that the agent has committed to in order to achieve some event-goals. A plan is appropriate if it is designed for the event in question (relevant) and its context condition is believed true (applicable). In contrast with traditional planning, execution happens at each step. The use of plans’ context-preconditions to make choices as late as possible, together with the built-in goal-failure mechanisms, ensures that the system is responsive to changes in the environment. In this paper we focus on the plan library to investigate ways of learning appropriate or better plan selection based on experience.

By grouping together plans responding to the same event (or goal) type, the plan library can be seen as a set of *goal-plan tree* templates: a goal node has children representing the relevant plans for achieving it; and a plan node, in turn, has children representing the subgoals (including primitive actions) of the plan. These structures can be seen as AND/OR trees: for a plan to succeed all subgoals and actions of the plan must be successful (AND); for a subgoal to succeed one of the plans to achieve it must succeed (OR).

Consider, for instance, the hierarchical structure shown in Figure 1(b). A link from a goal to a plan means that this plan is relevant (i.e., potentially suitable) for achieving the goal (e.g., $P_1 \dots P_4$ are the relevant plans for event goal G); whereas a link from a plan to a goal means that the plan needs to achieve that goal as part of its (sequential) execution (e.g., plan P_A needs to achieve goal G_{A1} first and then G_{A2}). For compactness, an edge with a label $\times n$ states that there are n edges of such type. Leaf plans directly interact with the environment and so, in a given world state, they can either succeed or fail when executed; this is marked accordingly in the figure *for some particular world* (of course, in other states, such plans may behave differently). In some world, given successful completion of G_A first, the agent may achieve goal G_B by selecting and executing P_B , followed by selecting and executing two working leaf plans to resolve goals G_{B1} and G_{B2} . If the agent succeeds with goals G_{B1} and G_{B2} , then it succeeds for plan P_B , achieving thus goal G_B and the top-level goal G itself. There is no possible

successful execution though, if the agent decides to carry on any of the three plans labelled P'_{B2} for achieving the low-level goal G_{B2} .

As can be seen, *plan-selection* is critically important. Standard BDI systems leverage domain expertise by means of the context conditions of plans. In this work, we are interested in exploring how a situated agent may *learn* plan selection based on experience, in order to improve goal achievement.

2.2. Learning for BDI Plan Selection

In order to facilitate learning regarding which plan should be executed for a given goal in a particular world state, we first replace each plan's boolean formula that is the standard representation for context conditions in BDI programming languages, with a decision tree [13] that provides a judgement as to whether the plan is likely to succeed or fail for the given situation.

To select plans based on information in the decision trees, we use a probabilistic method that chooses a plan based on its believed likelihood of success in the given situation. This approach provides a balance between exploitation (we choose plans with relatively higher success expectations more often), and exploration (we sometimes choose plans with lower success expectation to get better confidence in their believed applicability by trying them in more situations). This balance is important because ongoing learning influences future plan selection, and subsequently whether a good solution is learnt.

The resolution of goals in BDI execution results in the invocation of plans that in turn may post sub-goals that are further handled by sub-plans in a hierarchical manner. In a programming context, this is equivalent to making a function call that in turn calls sub-functions. However a sub-goal may have a number of possible plans for achieving it, some of which will work better in particular situations than others. In our learning context, where we do not yet know which plans work well in which situations, a plan may fail not because the plan was a bad choice in the given situation, but instead because the run-time choice of sub-plans was incorrect for the situation.

Our first approach [14] to address the learning problem was that of careful consideration whereby failures are recorded for learning purposes only when we are sufficiently sure that the failure was not due to poor sub-plan choices. We have shown that this conservative approach is more robust, though often slower, than a more aggressive approach which records all experiences, but can in some particular cases completely fail to learn.

Our second approach reported in [15] was to adjust the plan selection probability based on some measure of our confidence in the decision tree. We consider the reliability of a plan's decision tree in a given world state

to be proportional to the number of sub-plan choices (or paths below the plan in the goal-plan hierarchy) that have been covered in that world state. Here *coverage* [15] refers to the set of explored paths relative to the set of all possible paths. The greater the coverage, the more we have explored and the greater the confidence in the resulting decision tree. By biasing the plan selection probability with a coverage-based confidence measure we achieved the same robustness as that of conservative recording of failure cases. The coverage approach, however, is more flexible as the extent to which this is used can be readily adjusted by parameters in the selection formula.

A limitation with the previous approaches is that events were assumed to be propositional atoms, i.e., parameterised event-goals were not considered. By *parameterised* we mean an event-goal that may contain “data” as part of its definition. For instance, event *travelTo(dest)* may represent the goal to travel to location *dest*. In general, a goal to move to location *A* may require different strategies than those for addressing a goal to move to location *B*, and the learning must account for this. Another limitation is the assumption that the agent’s plan library does not include recursive subgoaling, so that the goal-plan tree structure induced is always *finite*. For example, the above plan rule in the elevator controller’s library for handling subgoal *GoTo(floor)* would not be allowed, since its procedure involves posting the same subgoal event as the rule’s head. Clearly both limitations would preclude the applicability of the approach in many practical domains where hierarchies are usually expressed in a compact manner by using parameterised goal events and plans, and often make use of (direct or indirect) recursive procedures to encode iterative strategies.

Furthermore, the coverage approach [15] does not scale to recursive structures. Conceptually we can unfold the recursive structure to a specified depth. However, the number of paths is exponential in the recursion number and further compounded by parameterised event-goals and the number of possible world states. An additional limitation is that coverage does not consider domain complexity. For instance, a leaf plan that has no sub-goals will achieve full coverage when it is tried once, after which selection will be fully biased towards the plan’s decision tree classification. However, the decision tree that at this point has only witnessed one world will generalise the outcome to all as-yet-unseen worlds leading to misclassification.

In the following section we present the details of our learning approach, incorporating both parametrised goals, and recursion, as well as a new simpler confidence measure that is based on the general idea of coverage but does not suffer from its limitations.

3. The BDI Learning Framework

Our learning task is as follows: *Given past execution data and the current world state, determine which plan to execute next in order to best address the given event-goal.* In the BDI sense, our task is to learn the context condition of each plan in the goal-plan hierarchy. In this section we describe our BDI Learning Framework that enables such learning. In particular we describe the use of decision trees for learning context conditions and the confidence-based probabilistic plan selection that incorporates this learning, while focusing on parameterised event-goals and recursion.

3.1. Integrating Decision Trees into Context Conditions for Plans

A plans context condition is a logical formula that is constructed at design time and evaluated against an event-goal at run time to determine if the plan is applicable in the given world state.² To allow the context condition to be learnt over time, we annotate each plan’s context formula with a *decision tree*.³[14] The idea is that the agent starts with some *necessary but possibly insufficient* conditions for each plan (provided by the designer), and over time and in the course of trying plans in various world states will *refine* each plan’s context condition using the learnt decision tree.

The choice of decision trees for learning is motivated by several factors. Firstly, decision trees support hypotheses that are a disjunction of conjunctive terms and this representation is compatible with how context formulas are generally written. Secondly, decision trees can be converted to *if-then* rules that are human readable and can be validated by a domain expert. Finally, decision trees are robust against training data that may contain errors. This is specially relevant in stochastic domains where applicable plans may nevertheless fail due to unforeseen circumstances.

For each plan, the training set for its decision tree contains samples of the form $[w, o]$, where w is the world state in which the plan was executed and o was the boolean outcome (success or failure). The world state w itself is a set of discrete attributes that together represent the state of affairs. Initially the training set is empty and grows as the agent tries the plan in various world states and records each result. Over time the decision tree learnt from

²Context formulas may reference internal beliefs as well as environment states, and for this study we treat both as included in the world state.

³It is perfectly feasible to combine the existing logical formula with the decision tree classification, but to aid our understanding of the decision tree learning in this study we always use an empty initial formula.

the training set will contain only those attributes of world state w that are relevant to that plans context condition.

The number of attributes in world state w and their range has a bearing on the size of the training set required to correctly learn the context condition. In general, world state w should be constructed with all attributes that are possibly relevant to the context condition. For instance, for a plan to pick objects using a robotic arm, the attributes *objectSurface* and *gripperWet* are likely relevant and should be included, while the attribute *dayOfWeek* possibly is not and may be excluded. The choice of attributes to include in the world state w is eventually a design decision and dependent on domain knowledge. For our purposes we assume that the designer provides a set of all attributes that are considered *possibly relevant* to the context condition of the plan.⁴ In the worst case, this set is the full set of available attributes.

The decision tree inductive bias is a preference for smaller trees. In other words, the induction will trade-off some accuracy in classification for compactness of representation. This means that some training samples get incorrectly classified in the wrong “bucket” (where the bucket name is *success* or *failure* in our case) when the actual outcome class of those training samples is different. So in a given bucket (or class) one may get a total of m training samples out of which n are misclassified. The ratio $1 - (n/m)$ gives the likelihood of class membership,⁵ and is the fraction we use as the expected likelihood of success of the plan.

3.2. Support for Recursive Event-Goals

Recursion in our context refers to the case where the resolution of an event-goal instance $G(\vec{x}_1)$ involves first the resolution of goal-event instance $G(\vec{x}_2)$ of the same type G . The result is a growing stack of pending $G(\vec{x}_i)$ event-goals that eventually terminate in $G(\vec{x}_n)$ whose parameters satisfy the termination conditions where a non-recursive plan choice is made.

In order to understand the impact of recursion on context learning, we use the notion of an *execution trace* of the form $G_0(\vec{x}_0)[P_0 : w_0] \cdot G_1(\vec{x}_1)[P_1 : w_1] \cdot \dots \cdot G_n(\vec{x}_n)[P_n : w_n]$, that represents a sequence of event-goals along with the plans selected to handle them and the world state in which the selections

⁴An automatic compilation of potential relevant propositions can be done by analysing, if available, the preconditions and effects of actions that might be executed when handling a goal. This, however, is out of the scope of this paper.

⁵In our study we use algorithm **J48**, a version of **c4.5** [13], from the **weka** learning package [17] that automatically provides this ratio.

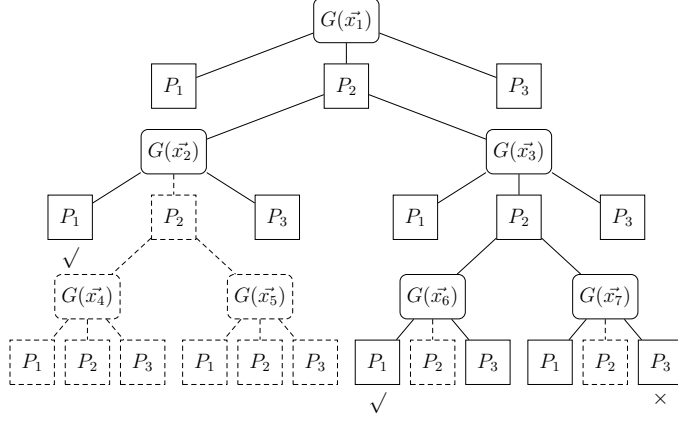


Figure 2: Goal-plan hierarchy containing a parameterised goal G handled by three plans P_1 , P_2 and P_3 . Here plan P_2 posts two instances of G resulting in recursion. Two levels of recursive unfolding are shown. Dashed P_2 nodes indicate unexplored recursive sub-trees.

were made. So $G_i(\vec{x}_i)[P_i : w_i]$ captures the case where plan P_i was selected in world state w_i in order to achieve the goal-event $G_i(\vec{x}_i)$.

Consider the example BDI goal-plan hierarchy of Figure 2. The structure has just a single parameterised goal G and three options to handle it, one of which, P_2 , in turn posts two instances of the same parameterised goal G . In this way, the only plans that take an action in the environment are P_1 and P_3 . The figure highlights an execution trace as follows:

$$\lambda = G(\vec{x}_1)[P_2 : w_1] \cdot G(\vec{x}_2)[P_1 : w_1] \cdot G(\vec{x}_3)[P_2 : w_2] \cdot G(\vec{x}_6)[P_1 : w_2] \cdot G(\vec{x}_7)[P_3 : w_3].$$

The first choice in the execution results in the selection of plan P_2 to handle event-goal instance $G(\vec{x}_1)$ in a given world w_1 . Plan P_2 in turn immediately posts the event-goal instance $G(\vec{x}_2)$ that is successfully handled by the non-recursive node P_1 . Plan P_2 then posts the second event-goal instance $G(\vec{x}_3)$, which then is handled by itself in a recursive manner. The outcome is that λ traces a path that involves the successive execution of leaf plan P_1 for event-goal $G(\vec{x}_2)$ followed by another execution of P_1 for event-goal $G(\vec{x}_6)$, finally terminating in the failure of leaf plan P_3 for event-goal $G(\vec{x}_7)$. Note that if plan P_2 had instead been selected to handle $G(\vec{x}_7)$ then a deeper recursive call would have ensued. Similarly if earlier in the execution trace plan P_2 was selected to handle event-goal $G(\vec{x}_2)$ then a different recursive sub-tree (shown in Figure 2 as dotted nodes under $G(\vec{x}_2)$) would have unfolded.

The immediate implication of a recursive goal-plan structure is that the

size of the hierarchy is no longer static but instead unfolds in a dynamic manner. The risk then is that since the conditions that terminate recursion are not ready at the start (we are trying to learn them), then the agent may get trapped in an infinite recursive loop during exploration. This has implications for any strategy that relies on the structure being finite. For instance, our conservative recording approach [14] and coverage-based confidence measure [15] both suffer from this problem. Incidentally, the simpler aggressive recording approach [14] is not impacted by recursion as it does not consider the goal-plan structure.

One way to resolve this issue is to treat all recursive goals simply as sub-trees in a static structure and limit the recursive *unfolding* to a maximum allowed depth. In this study we use this bounded recursion approach for handling recursive structures. It follows then that wherever a recursive structure applies, a maximum recursion value must be supplied. This may not be an unrealistic requirement given that the domain expert will usually have some idea of how much recursion is sufficient for a given parameterised event-goal.

3.3. Calculating Confidence in the Decision Tree Classification

The typical use of decision trees lies in the *offline* induction from a complete training set. In that sense, the use of decision trees in our framework is unorthodox since the training set is built incrementally by recording samples after each new execution. This results in the training set being incomplete in the early stages of learning,⁶ leading to misclassification. A confidence measure in the decision tree classification is therefore desirable to address this issue. Previously [15] we showed how the coverage of possible execution paths below the plan in the goal-plan hierarchy may be used to build such a measure. Here we propose a new confidence measure that builds on the idea but that does not suffer from its limitations (Section 2.2).

Our requirement for the confidence measure is that it be a monotonic function whose values transition from no confidence (0.0) to full confidence (1.0) based on experience. Specifically, the experiences we are interested in should constitute coverage of the plan complexity (number of sub-plan choices) and the domain complexity (number of world states in which the plan applies). Since an exact calculation of such coverage does not scale for all practical purposes then we are interested in an *approximate* coverage that is still representative of the state of affairs but is simpler to compute.

⁶Training data is incomplete in the sense that the agent has only collected a portion of the full data set required to learn the correct classification.

One way to achieve this is to use a monotonic decay function⁷ (for instance $\epsilon_i = \epsilon_{i-1} * \delta$ where $\delta < 1.0$) but where the decay factor δ is tied to the complexity involved. This way, a plan that has a larger number of sub-plan choices will utilise a slower decay factor δ taking longer to reach full confidence $(1 - \epsilon)$ than another plan that has less choices to make. For goal-plan complexity this decay δ_{Pt} may be calculated offline by analysing the goal-plan hierarchy. In this work, we have calculated δ_{Pt} in terms of average breadth and depth of the structure, where depth is the maximum level of recursion in this case, to provide an approximation of the complexity of the structure.⁸ A similar treatment is possible for domain complexity although the decay factor in this case cannot be pre-determined since the number of world states is not known upfront and is dependent on the domain. For domain complexity then, it may be reasonable to treat the decay factor δ_{Pd} as a parameter specified by the domain expert.

$$c_P = (1.0 - \epsilon_{Pt}) * (1.0 - \epsilon_{Pd}). \quad (1)$$

Equation 1 shows how the final confidence c_P is calculated for a given plan P . Here ϵ_{Pt} is the plans tree complexity decay while ϵ_{Pd} is the plans domain complexity decay. The actual updates to the decay values are performed each time the plan P is executed while the rate of decay is governed by the decay factors δ_{Pt} and δ_{Pd} accordingly.

3.4. Handling Parameterised Event-Goals

Our BDI learning framework account presented earlier [14, 15] did not account for *parameterised event-goals* but only for event-goal instances. In practical BDI systems, it is often the case that a single plan will handle all instances of a parameterised event-goal. Furthermore event-goal instance parameters will generally be included in the context logical formula.

Consider again the goal-plan structure in Figure 2 and the highlighted solution path terminating in the leaf plans indicated by the \checkmark symbol. An important point here is that the indicated solution applies to the event-goal *instance* $G(\vec{x}_1)$ and to that instance alone. For a different instance $G(\vec{y}_1)$

⁷This technique is frequently applied in machine learning algorithms for balancing between exploration of choices and exploitation of learning.

⁸There are other ways of calculating δ_{Pt} (e.g., in [15] we have used an accurate calculation of the number of choices below plans; however this is not feasible anymore when goals and plans are schemas with possibly infinite instances); the main idea is to somehow measure the complexity of a hierarchical structure, and is the subject of ongoing work.

the solution path would likely be different (one way to visualise this in the Figure 2 is to think of it as an animation where the event-goal parameters and the placement of the \surd symbols changes on each frame). This means that event-goal instance parameters must also be considered as input for a plan’s decision tree in order to learn solutions per event-goal instance.

We include such an account by augmenting the training samples for the decision tree with the event-goal parameters. As such, the training set now contains samples of the form $[w, x, o]$ where the world state w is the initial set of all relevant attributes that represent the state of affairs, x is the set of all event-goal parameters, and o is the outcome class (success or failure). Incorporating the event-goal parameter set x in the training data is sufficient for learning with parameterised event-goals, and no fundamental change to the framework is required.

3.5. Calculating Plan Selection Weights based on Confidence

Typical BDI platforms offer several mechanisms for plan selection from a set of applicable plans, such as plan precedence and meta-level reasoning. However, since these mechanisms are pre-programmed and do not take into account the experience of the agent, we provide a new *probabilistic plan selection* function for this purpose.

For each plan, given its expectation of success (as determined by its decision tree learning) and a confidence measure in this expectation (based on coverage), we calculate a final *selection weight* that is indicative of the likelihood of the plan being selected for execution. Equation 2 shows how the plan selection weight $\Omega_P(w)$ is calculated for a given world state w .⁹ Initially, the confidence c_P is zero and the weight takes the default value of 0.5. Over time, as the confidence improves towards the final value of 1.0, the selection weight approaches the value $\kappa_P(w)$ estimated by the plan’s decision tree.

$$\Omega_P(w) = 0.5 + [c_P * (\kappa_P(w) - 0.5)]. \quad (2)$$

Given the set of applicable plans for resolving event-goal G in world state w then, our probabilistic plan selection mechanism chooses a plan P_i with a probability directly proportional to its selection weight $\Omega_{P_i}(w)$. Such selection ensures a balance between the *exploitation* of current know-how and the *exploration* of new choices that is necessary for online learning tasks.

⁹The formulation of the plan selection weight is described in [15].

4. A Case Example: The Hanoi Towers Robot

To evaluate our learning framework we consider an existing BDI program from the JACK agent platform distribution [10]. The example involves a robot playing the well-known Towers of Hanoi game where the goal is to stack discs of decreasing size onto a single pin. The rules of the game forbid discs to be moved onto smaller discs, however top discs may be moved onto discs of larger size across three pins. The problem is interesting for our purposes since the example solution makes use of parameterised event-goals and recursion. Furthermore, unlike our previous evaluations [14, 15] with synthetic plan libraries, here the evaluation criteria is clear: *does our learning framework achieve the performance of the existing system?*

The example solution consists of a *Player* agent that solves the game for any given legal initial configuration. The game solving strategy is encoded in plan *DiscStacker* that solves for one disc at a time starting from the largest and ending with the smallest onto a chosen pin. This in turn is achieved by posting event $Solve(d,p)$ for solving disc d onto pin p . There are four plans that are relevant for this purpose:

SolveRight This plan solves moving a disc to the pin it is already on. Since the goal is already true, the plan does *nothing*.

SolveTopMove This plan moves the disc d to the destination pin p if the disc is not already there and if the move is legal. The actual move is performed by the primitive action $move(p2,p)$, where $p2$ is the source pin of disc d .

SolveTop This plan solves for the case when the disc d may be legally lifted but cannot be legally placed at the destination because the top disc on the destination pin is smaller than d . In this case, the plan first moves all the discs in the destination pin that are smaller than disc d to the third (auxiliary) pin, and then re-posts the sub-goal to move d to pin p i.e. $Solve(d,p)$.

SolveMiddle This plan solves moving a disc from the *middle* of a stack. In this case, the plan first clears the source pin so that disc d becomes the new top of the pin. This is done by solving for sub-goal $Solve(d2,p2)$ where disc $d2$ is the disc currently on top of d and $p2$ is the auxiliary. Subsequently the plan re-posts the sub-goal of moving d to pin p i.e. event $Solve(d,p)$.

Figure 3 illustrates the goal-plan hierarchy for the domain. Here we focus on learning the recursive parameterised $Solve(d,p)$ event for which we remove the context conditions from the example plans and apply our framework.

4.1. Experimental Setup

The aim of this study is to evaluate our learning framework for recursive event-goals. For this reason our experimentation with the Hanoi problem

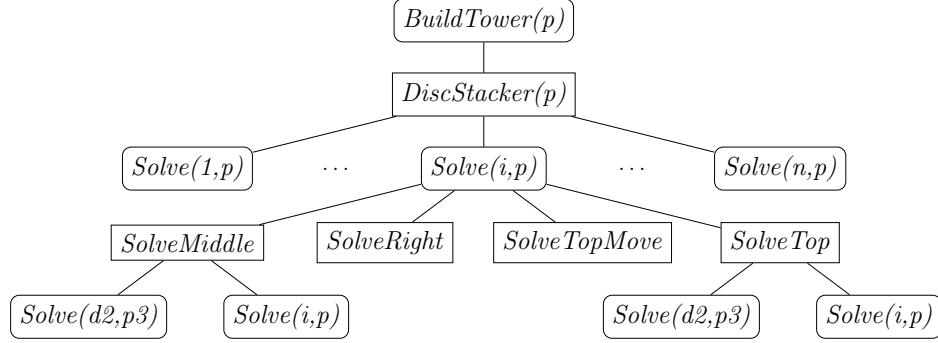


Figure 3: Goal-plan hierarchy for the Towers of Hanoi domain.

focuses on learning to resolve the recursive event *Solve* only and not on learning the strategy that solves the full Hanoi towers problem (this is done by *DiscStacker(p)*). Since the full set of possible *Solve* events and initial pin configurations is large, our first step is to construct a sufficiently rich subset that we will use to evaluate our learning approaches. We proceed by running the original Hanoi program for a number of randomly generated *Solve* events. For each run we record the *Solve* event, the initial pins configuration, and the maximum recursion encountered for the solution. This gives us a bag of several initial configurations for each recursion level that is a subset of all possible configurations.

Next, we run each candidate approach on the set of saved configurations for a given recursion level. i.e. where all solutions lie exactly at the specified recursion number. We use a fixed random generation seed for each experiment so that the same sequence of *Solve* events is generated for each learning approach. This isolates any environmental factors and allows us to attribute any differences in performance to the learning approaches alone.

4.2. Results

The following results are for a Hanoi problem with *five* discs.¹⁰ Each plans domain complexity decay factor for the confidence calculation of Equation 1 are set to $\delta_{Pd} = 0.9$. For goal-plan tree complexity, we use $\delta_{Pt} = 1$ for non-recursive plans and $\delta_{Pt} = [1 - (1/r^k)]$ for recursive plans. Here r is

¹⁰We use five discs in order to keep the state space rich enough yet sufficiently small to allow learning runs to be completed and evaluated in reasonable time.

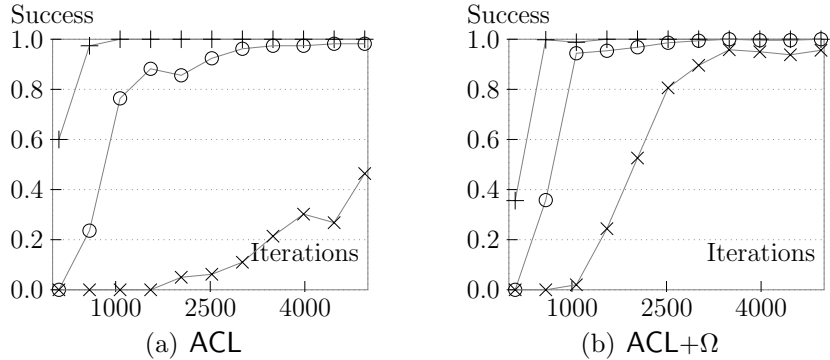


Figure 4: Agent performance under ACL and ACL+ Ω schemes for solutions at recursion levels one (pluses), three (circles) and five (crosses). Each point represents results from 5 experiment runs using an averaging window of 100 samples.

the run-time recursion level and k is arbitrarily set to 4.0,¹¹ and controls the rate of change of decay i.e. the number of steps to reach full confidence. In all experiments, the recursion is bound to a maximum of *eight* levels that is sufficient to solve all configurations for a five-disc Hanoi problem. The performance of two learning configurations is contrasted. The baseline learning algorithm ACL refers to the original aggressive learning approach of [14] and [15] using the original probabilistic plan selection function that has no confidence-based bias (uses decision tree expectation of success only). The new algorithm is referred to as ACL+ Ω and uses the same aggressive learning approach as the former but combined with the new confidence-based probabilistic selection function (Equation 2) presented in this study.

Experiment 1. To understand how the two approaches perform for solutions of varying difficulty we conducted a set of tests with solutions at different recursive levels. Each test consisted of resolving a known set of *Solve* event configurations, saved earlier as described in Section 4.1, whose solutions all required a given recursive depth. Figure 4 shows that as the solution difficulty increases from one to five recursion levels, ACL performance drops much more significantly compared to that of ACL+ Ω . For instance, for solutions requiring five levels of recursion (crosses in Figure 4), ACL achieves only 50% success at 5k iterations whereas ACL+ Ω achieves 95% success by 3.5k iterations. The poor performance of ACL may be attributed to the fact that

¹¹In future work we hope to establish principles for determining general parameters.

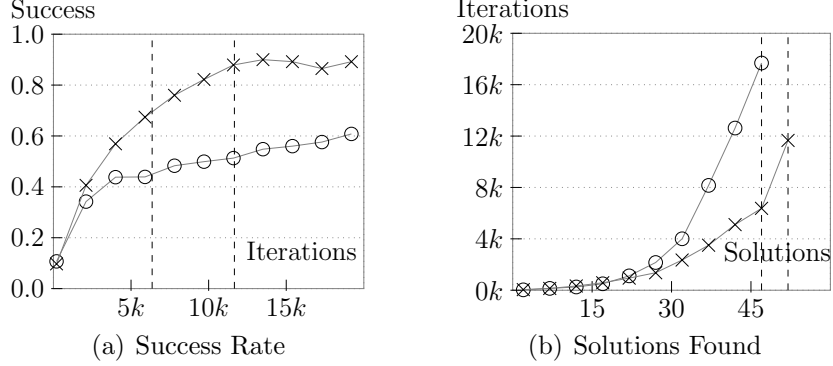


Figure 5: Agent performance under **ACL** (circles) and **ACL+Ω** (crosses) schemes. Each point represents an average result from 5 experiment runs.

deeper solutions require more $move(p2, p)$ steps and where an earlier success does not exist to guide selection at each resulting state the exploration is mostly random. On the other hand, the confidence-based measure of Equation 2 takes into account the goal-plan tree complexity and is able to guide the **ACL+Ω** exploration towards the deeper solutions.

Experiment 2. Next, we conducted an experiment that consisted of resolving the full set of saved *Solve* events i.e. the set of all solutions for recursion levels one to five. Figure 5(a) shows the results for the two approaches where **ACL+Ω** performs better than **ACL** as expected from the previous experiment results. For the same experiment, we also recorded the number of solutions found. Figure 5(b) shows that **ACL+Ω** resolves all 52 goals within 12k iterations whereas **ACL** resolves only 47 by the end of the experiment at 20k. At a similar point of comparison, to resolve 47 goals **ACL+Ω** takes around 6.4k iterations whereas **ACL** takes more than twice as long at 17.7k. The vertical dashed lines in Figure 5(a) and Figure 5(b) mark the 47th and 52nd solutions.

Experiment 3. Finally, we ran a third experiment to understand the impact of applicability thresholds. In the classical BDI framework plans are either applicable or not in a boolean decision. However, in our modified framework plans are applicable according to the selection weight given by Equation 2. Since plan execution is often not cost-free in real systems, it is likely that an adequate plan selection scheme would not select *any* plan if they all have too low an expectation of success, and an agent may fail a plan without even

trying. To represent this scenario we setup an applicability threshold of 20% whereby plans with expectations of success less than this threshold would not be selected. In this case ACL shows a complete inability to learn as reported earlier [15]. In contrast ACL+ Ω benefits from the adjusted selection weights of Equation 2 and shows similar performance as before (Figure 5(a)).¹²

Analysis. Observe in Figure 5(a) that ACL+ Ω does not reach the performance of the hand-crafted JACK program and converges to about 90% success even though it successfully discovers all solutions (Figure 5(b)). This is because the decision tree representation does not guarantee that the training data will always be correctly classified (we discuss this accuracy versus compactness trade-off in Section 3.1). For instance, a decision tree may report a poor likelihood of success for a given state even when the associated training sample indicates success, due to the sample being misclassified in the failure “bucket”. One could guarantee correctness by referencing the training data directly, for instance using a look-up table. However, the decision tree representation is preferable for its compact representation combined with the ability to generalise to as-yet-unseen world states. Currently, when our learnt decision trees are converted to rules they do not “look” like the original ones but are far more complex. This is mainly due to representational differences as our simple representation is propositional whereas the original conditions are relational.¹³

5. Discussion and Conclusion

This paper builds on earlier work [14, 15] that extends the typical BDI programming framework to use decision trees as (part of) a plan’s context condition, with a probabilistic plan selection mechanism that caters for both exploration and exploitation of plans. Previously we have shown that due to the structure of BDI programs, care must be taken in how learning is used, to avoid problems in certain situations. In some cases these problems lead to failure to learn at all, as we also show here.

In this paper we extend previous work to allow for parameterised goals (e.g. *travelTo(dest)*) and also for recursion, both of which are necessary for

¹²The threshold of 0.2 while somewhat arbitrary is consistent with that used earlier [15]. The difference between the default weight (0.5) and the threshold weight (0.2 in this case) decides how much “give” we have in the exploration. The closer the threshold is to 0.5 the greater the chance that the plan will be aborted before a solution is found.

¹³We hope to address this using relational decision trees in future work.

real applications. In doing this our previous confidence measure which relied on a finite goal-plan tree did not scale, so we have provided an approximate measure, relying on the principles that have been shown correct, but without the limitations. This paper also takes an existing BDI program involving parameterised goals and recursion, and evaluates our approach using this program. By removing the existing context conditions, and then learning the correct behaviour, we show that we are able to obtain good (although not perfect)¹⁴ performance. We also demonstrate that the naive approach to learning, that does not account for the BDI program structure fails to learn given some program structures and an applicability threshold.

There is still work to be done before our framework may be applied for *practical* on-line learning in situated agents. Firstly, the framework has not been integrated with standard BDI failure handling and recovery. Clearly this will be needed (and is the subject of ongoing work), but we do not expect this to undermine any results described here. In fact a careful integration of failure handling could improve the speed of learning as multiple attempts could be made to achieve a (sub)-goal. However care needs to be taken regarding changes to world state and possible interactions between failed attempts and eventually successful ones.

Secondly, our use of decision trees is naive. For instance, currently execution data is maintained forever and decision trees re-built after each plan execution. Furthermore, we learn using actual world states, whereas an improvement would be to learn using relational world information. While not ideal, our setup nonetheless allows us to focus on the nuances of learning in BDI programs first without worrying about the underlying techniques.

Finally, we do not detect and learn interactions between sibling goals in the context of a particular parent; each subgoal is treated “locally.” To handle such interactions, the selection of a plan for resolving a sub-goal should also be predicated on the goals higher than the sub-goal, that is, it should take into account the “reasons” for the sub-goal. Addressing this would require substantial modification to the BDI programming style in terms of representation, which is out of the scope of this work.

The issue of combining learning and deliberative approaches for online decision making in BDI-like systems has not been widely addressed. Hernández et al. [18] give a preliminary account of how decision trees may be induced on plan failures in order to find alternative logical context conditions in a

¹⁴We would hope that when learning is combined with programmer provided context conditions, the problems preventing perfect learning here would be avoided.

deterministic paint-world example. In [19] learnt user preferences are incorporated during BDI plan selection in a dialogue manager application using a decision tree learner. In contrast, [20] take the approach of refining existing BDI plans or learning new plans as a sequence of recorded actions based on prescriptions provided by the domain expert. In [21], low level robot soccer skills are learnt offline and then used in the deliberative decision making process once deployed. More recently, [22] give a comprehensive account of integrating learning in BDI deliberation for a real world ship berthing logistics domain. Here a neural network module is first trained offline on the available shipping port data and then used in a deployed BDI system to improve plan selection. Their results show significant improvement in berth productivity over the existing system of human operators.

A closely related area to BDI is that of hierarchical task network (HTN) systems where task decompositions used are similar to BDI goal-plan hierarchies. Recently, in similarly motivated work to ours, [23] proposed a method for learning HTN method preconditions under partial observations. There, a set of constraints are constructed from observed decomposition trees that are then solved *offline* using a constraint solver. In contrast, in our work learning and deliberation are fully integrated in a way that one impacts the other and the classical exploration/exploitation dilemma applies.

The BDI architecture has also been shown [24] to be related to Markov Decision Processes that are heavily used for solving optimisation problems in reinforcement learning [13]. A sub-area of work related to ours is hierarchical reinforcement learning [25] where task hierarchies similar to BDI are used. When the aim is to find locally optimal solutions for each sub-MDP in the hierarchy, similar issues as ours arise, such as goal inter-dependence. In general, global optimality is possible only when information is fed into the sub-task (i.e. value functions use the entire state space), consistent with our analysis of goal inter-dependence issues. Interestingly, work by Dietterich [26] also supports the use of simultaneous learning at all levels (similar to our ACL based approaches) instead of waiting for the children to converge (analogous to our conservative approach [15]).

Although there is still work to do before we can expect learning to be successfully integrated into a fully autonomous BDI agent, the work reported here is significant in that it provides a solid foundation for adding new capabilities to BDI agents to allow them to learn and adapt based on experience.

- [1] N. R. Jennings, An agent-based approach for building complex software systems, Communications of the ACM 44 (4) (2001) 35–41.

- [2] R. A. Belecianu, S. Munroe, M. Luck, T. Payne, T. Miller, P. McBurney, M. Pechoucek, Commercial applications of agents: Lessons, experiences and challenges, in: *Proceedings of Autonomous Agents and Multi-Agent Systems (AAMAS)*, ACM Press, Hakodate, Japan, 2006, pp. 1549–1555.
- [3] M. Ljungberg, A. Lucas, The OASIS air-traffic management system, in: *Proceedings of the Pacific Rim International Conference on Artificial Intelligence (PRICAI)*, Seoul, Korea, 1992.
- [4] Z. Tu, M. Ferry, G. Prickett, C. Heinze, Truly autonomous UAVs and teaming, in: *Twelfth Australian Aeronautical Conference*, 2007, pp. 483–489, engineers Australia.
- [5] M. E. Pollack, The uses of plans, *Artificial Intelligence Journal* 57 (1) (1992) 43–68.
- [6] M. Bratman, D. Israel, M. Pollack, Plans and resource-bounded practical reasoning, *Computational Intelligence* 4 (4) (1988) 349–355.
- [7] M. Bratman, *Intentions, Plans, and Practical Reason*, Harvard University Press, 1987.
- [8] D. C. Dennet, *The intentional Stance*, MIT Press, 1987.
- [9] S. S. Benfield, J. Hendrickson, D. Galanti, Making a strong business case for multiagent technology, in: *Proceedings of Autonomous Agents and Multi-Agent Systems (AAMAS)*, ACM Press, 2006, pp. 10–15.
- [10] P. Busetta, R. Rönquist, A. Hodgson, A. Lucas, JACK intelligent agents: Components for intelligent agents in Java, *AgentLink Newsletter* 2 (1999) 2–5, Agent Oriented Software.
- [11] A. Pokahr, L. Braubach, W. Lamersdorf, JADEX: Implementing a BDI-infrastructure for JADE agents, *EXP - in search of innovation (Special Issue on JADE)* 3 (3) (2003) 76–85.
- [12] R. Bordini, J. Hübner, M. Wooldridge, *Programming Multi-agent Systems in AgentSpeak Using Jason*, Wiley Series in Agent Technology, Wiley, 2007.
- [13] T. Mitchell, *Machine Learning*, McGraw Hill, 1997.
- [14] S. Airiau, L. Padham, S. Sardina, S. Sen, Enhancing adaptation in BDI agents using learning techniques, *International Journal of Agent Technologies and Systems (IJATS)* 1 (2) (2009) 1–18.

- [15] D. Singh, S. Sardina, L. Padgham, S. Airiau, Learning context conditions for BDI plan selection, in: Proceedings of Autonomous Agents and Multi-Agent Systems (AAMAS), Toronto, Canada, 2010, to appear.
- [16] M. P. Georgeff, F. F. Ingrand, Decision making in an embedded reasoning system, in: Proceedings of IJCAI, Detroit, USA, 1989, pp. 972–978.
- [17] I. Witten, E. Frank, Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations, Morgan Kaufmann, 1999.
- [18] A. Guerra-Hernández, A. E. Fallah-Seghrouchni, H. Soldano, Learning in BDI Multi-agent Systems, Vol. 3259 of LNCS, Springer, 2004, pp. 218–233.
- [19] A. Nguyen, W. Wobcke, An adaptive plan-based dialogue agent: integrating learning into a BDI architecture, in: Proceedings of Autonomous Agents and Multi-Agent Systems (AAMAS), ACM Press, New York, NY, USA, 2006, pp. 786–788. doi:10.1145/1160633.1160771.
- [20] S. Karim, B. Subagdja, L. Sonenberg, Plans as Products of Learning, IEEE Computer Society, Washington, DC, USA, 2006, pp. 139–145. doi:10.1109/IAT.2006.100.
- [21] M. Riedmiller, A. Merke, D. Meier, A. Hoffman, A. Sinner, O. Thate, R. Ehrmann, Karlsruhe brainstormers - a reinforcement learning approach to robotic soccer, in: RoboCup 2000: Robot Soccer World Cup IV, 2001.
- [22] P. Lokuge, D. Alahakoon, Improving the adaptability in automated vessel scheduling in container ports using intelligent software agents, European Journal of Operational Research 177 (3) (2007) 1985–2015.
- [23] H. Zhuo, D. Hu, C. Hogg, Q. Yang, H. Munoz-Avila, Learning HTN Method Preconditions and Action Models from partial Observations, in: Proceedings of the Twenty-First International Joint Conference on Artificial Intelligence (IJCAI-09), 2009.
- [24] G. I. Simari, S. Parsons, On the relationship between MDPs and the BDI architecture, in: Proceedings of Autonomous Agents and Multi-Agent Systems (AAMAS), ACM Press, New York, NY, USA, 2006, pp. 1041–1048. doi:10.1145/1160633.1160818.
- [25] A. Barto, S. Mahadevan, Recent Advances in Hierarchical Reinforcement Learning, Discrete Event Dynamic Systems 13 (4) (2003) 341–379.
- [26] T. Dietterich, Hierarchical Reinforcement Learning with the MAXQ Value Function Decomposition, Journal of Artificial Intelligence Research 13 (2000) 227–303.